

Gaussian Mixture Model

We assume that the dataset $D = \{x_i\}$ is generated by a probability distribution $P(x)$ that is just a weighted sum of K gaussians.

$$P(x) = \sum_{k=1}^K \pi_k N(x; \mu_k, \Sigma_k)$$

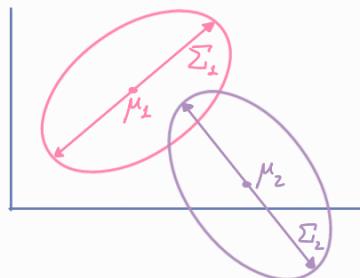
↑ prior probability ↑ mean ↑ covariance matrix

Total number of parameters of the model: $\underbrace{(\pi_k, \mu_k, \Sigma_k)}_{3} \times K$

Supposing K is known and that is an input for the model, we want to estimate the parameters.

Let's consider $K=2$, $\pi_1 = 0.6$

$$\pi_2 = 1 - \pi_1 = 0.4$$



Given the parameters of the distribution we can easily generate a dataset: first we sample the gaussians and then we keep all the points.

$$\left. \begin{array}{l} \pi_1, \mu_1, \Sigma_1 \\ \pi_2, \mu_2, \Sigma_2 \end{array} \right\} \longrightarrow D$$

We are interested into the opposite

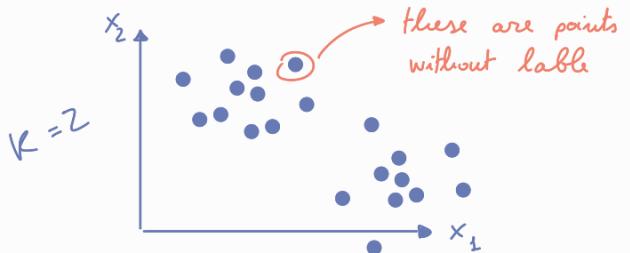
$$\left. \begin{array}{l} \pi_1, \mu_1, \Sigma_1 \\ \pi_2, \mu_2, \Sigma_2 \end{array} \right\} \xleftarrow{?} D$$

K-means

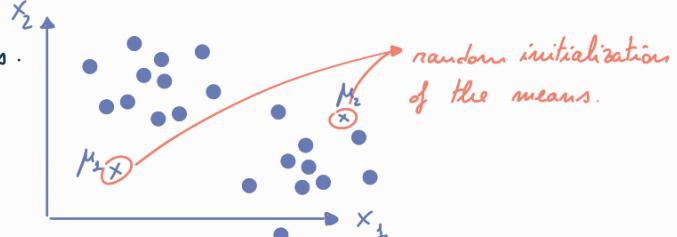
Simplify the problem: given a dataset $D = \{(x_n)\}$, given K , we want to compute the means for the K gaussians. This is called the K-means algorithm.

K-means is an iterative algorithm:

- ① The input is a set of unlabeled data and a value for K



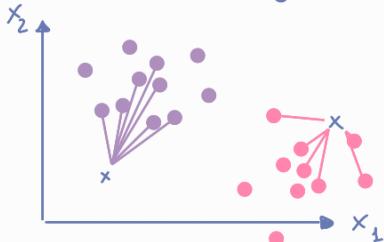
- ② Start with a random initialization of the means.



- ③ Form the clusters: assign each point in the dataset to the closest mean.

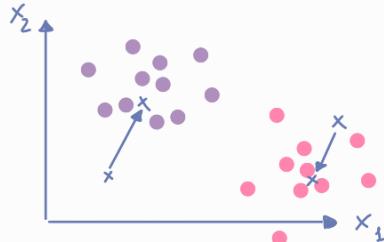
For each point we compute the distance from every mean and assign the point to the right cluster.

After this step we have partitioned the dataset in K subsets.



- ④ Refine the means given the clusters.

Compute the mean of each cluster, that will become the new estimation for the previous mean



- ⑤ Repeat the process until we have no changes between two iterations (clusters and means remains the same).

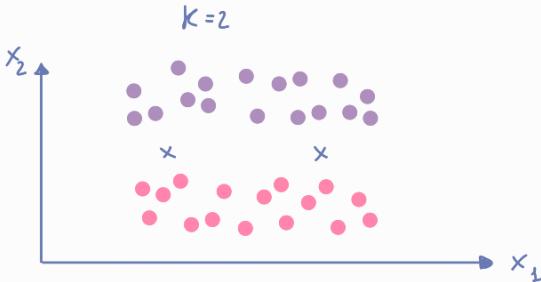
The convergence will always occur if the following two conditions are satisfied:

- ① Every time you iterate and form a new cluster with a new mean, the sum of distances from each training sample to that training sample's group centroid is decreased;
- ② There are only finitely many partitions of the training examples into K clusters.

Remarks on K-means

- The number of clusters K must be determined before hand;
- Sensitive to initial condition when few data is available;
- Not robust to outliers. Data very far from the real centroid may pull the mean far away;
- Does not account for the covariance matrix

Limitations mostly depends on the distance function



Even if you correctly guess K , since K-means is just considering a fixed covariance, with this dataset you will likely get this solution.

Improvements

- Use K-means clustering only if there is lot of data available;
- Use median instead of mean;
- Define better distance functions;

K-means still uses only a portion of the parameters of the model

Latent variables

Let's introduce the concept of **latent variables**, useful to describe the problem but not part of the dataset. The latent variables will express which gaussian has generated which point.

Introduce latent variable $z_k \in \{0, 1\}$ (boolean) with $z = (z_1, \dots, z_K)^T$, one for each gaussian, using a 1-out-of- K encoding: z_k is 1 if the point has been generated by gaussian k and all the other elements $z_1, \dots, z_{k-1}, z_{k+1}, \dots, z_K$ are 0.

$P(z_k = 1) = \pi_k$ → Probability that z_k is 1 corresponds to the prior probability π_k = probability of selecting gaussian k before looking at the data.

Since we assume that the samples are independently distributed, the probability of the vector z is just the product of the probabilities $\pi_k^{z_k}$

$$P(z) = \prod_{k=1}^K \pi_k^{z_k}$$

If we know z , we can express the conditional probability of x

$$P(x | z_k = 1) = N(x; \mu_k, \Sigma_k)$$

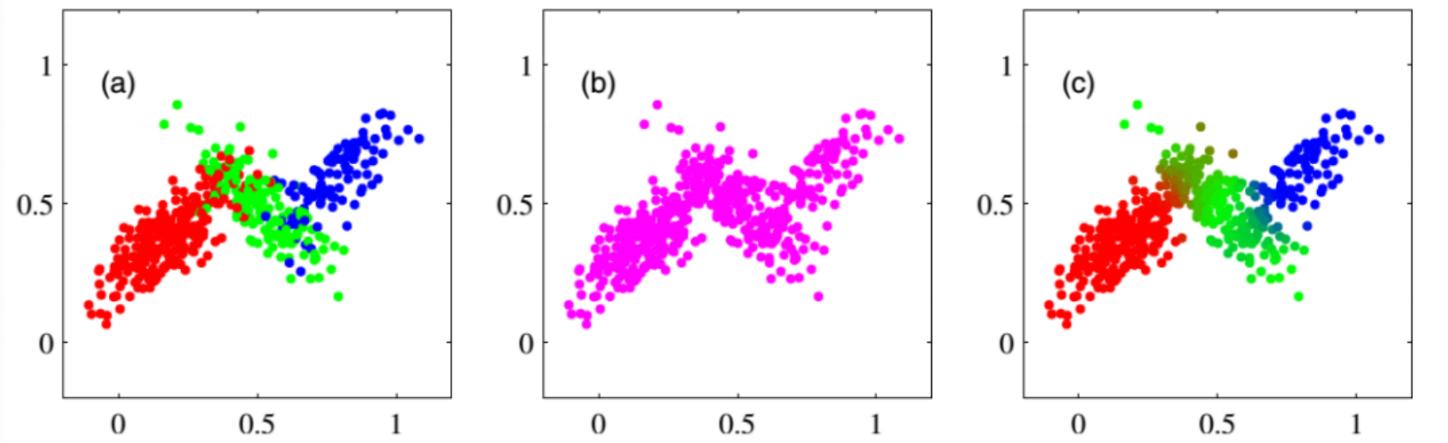
Thus (compact notation, putting together this notion for all the K gaussians) :

$$P(x | z) = \prod_{k=1}^K N(x; \mu_k, \Sigma_k)^{z_k}$$

We also know that (joint distribution) :

$$P(x, z) = P(x | z)P(z) \rightarrow \text{chain rule}$$

Joint probability distribution for a set of random variables gives the probability of every atomic joint event on those random variables.



- Ⓐ $P(x, z)$ with 3 latent variables z (red, green, blue)
 - Ⓑ $P(x)$ marginalized distribution
 - Ⓒ $P(z_{m,k})$ posterior distribution
- Ⓑ is the input to the model, Ⓢ is the result, Ⓠ is the ground truth.

Now what we can do is to express $P(x)$

$$P(x) = \sum_z P(z) P(x|z) = \sum_{k=1}^K \pi_k N(x; \mu_k, \Sigma_k)$$

↑
probability of x is, for all the gaussians, the probability that x is generated by that gaussian times the probability of that gaussian

$\sum_{k=1}^K \pi_k N(x; \mu_k, \Sigma_k)$ is the Gaussian Mixture Model

Our distribution $P(x)$ can be seen as the marginalization of a distribution $p(x, z)$ over variables z .

$$\sum_z P(z) p(x|z)$$

Marginalization: given a known joint distribution of two discrete random variables x and y , the marginal distribution of either variable - x for example - is the probability distribution of x when the values of y are not taken into account. This can be calculated by summing the joint probability distribution over all values of y .

$$p_x(x_i) = \sum_j p(x_i, y_j) \quad \text{and} \quad p_y(y_j) = \sum_i p(x_i, y_j)$$

$y \setminus x$	x_1	x_2	x_3	$p(y)$
y_1	$\frac{1}{32}$	$\frac{2}{32}$	$\frac{1}{32}$	$\frac{8}{32}$
y_2	$\frac{3}{32}$	$\frac{6}{32}$	$\frac{3}{32}$	$\frac{15}{32}$
y_3	$\frac{9}{32}$	0	0	$\frac{9}{32}$
$p_x(x)$	$\frac{16}{32}$	$\frac{9}{32}$	$\frac{6}{32}$	$\frac{32}{32}$

How can we use this result?

Given observations $D = \{(x_m)\}_{m=1}^N$, each data point x_m is associated to the corresponding variable z_m , which is unknown. If we had z_m this wouldn't be an unsupervised learning problem.

Let's define the posterior probability (probability of selecting a gaussian after looking the data):

$$p(z_k = 1|x) = \frac{p(z=1)p(x|z_k=1)}{p(x)} \quad \text{Bayes rule}$$

$$= \frac{\pi_k N(x; \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x; \mu_j, \Sigma_j)}$$